

文章编号:1005-3085(2010)02-0321-12

核实数据下单指标 EV 模型的经验似然推断*

刘 强¹, 薛留根²

(1- 首都经济贸易大学统计学院, 北京 100070; 2- 北京工业大学应用数理学院, 北京 100124)

摘 要: 对于单指标 EV 模型, 当解释变量带有测量误差, 尤其是替代变量和解释变量的关系不能确定时, 单指标模型的估计就会出现维数灾祸问题。本文提出了一种降维方法, 有效解决了上述问题, 并利用经验似然方法给出了未知参数的置信域。由于采用降维技术, 因而所得结果在生物工程、网络工程等领域的数据处理中具有较好的应用价值。

关键词: 单指标 EV 模型; 经验似然; 核实数据; 降维

分类号: AMS(2000) 62G05; 62J02

中图分类号: O212.7

文献标识码: A

1 引言

考虑如下单指标模型

$$Y = g(\beta_0^T X) + e, \quad (1)$$

其中 $g(\cdot)$ 为未知的一元联系函数, X 为 \mathbf{R}^p 上的随机向量, β_0 为 $p \times 1$ 的未知参数, 且有 $\|\beta_0\| = 1$, 这里 $\|\cdot\|$ 表示向量的 Euclidean 模。 e 为模型误差。在生物工程等领域中的数据处理中, 经常遇到所谓的“维数灾祸”问题, 由于单指标模型通过对 $\beta_0^T X$ 而不是 X 回归, 有效地避免了上述问题, 因而该模型在实际问题中具有广泛的应用前景。关于上述模型, 已有不少的文献讨论, Xue 和 Zhu^[1] 利用经验似然的方法给出了未知参数 β 的经验似然置信域。

在实际应用中, 解释变量 X 的精确值往往很难得到, 或者时间消耗太多, 或者代价太高, 这时直接观测的是带有测量误差的替代变量 \tilde{X} 而不是真实变量 X 。通常把变量带有测量误差的模型称为 EV (errors-in-variables) 模型。一般而言, 替代变量 \tilde{X} 和真实变量 X 之间的关系往往比较复杂, 例如 $\tilde{X} = \psi(X, \varepsilon)$, 其中 ε 为测量误差, $\psi(\cdot)$ 为一个任意的已知函数。在这种情况下, 要对未知参数 β_0 进行有效的推断往往是比较困难的, 然而利用核实数据可以很好地克服这一困难。基于替代数据和核实样本的推断已经引起了众多统计学者的重视。例如 Sepanski 和 Lee^[2] 研究了基于核实数据的非线性 EV 模型。Wang^[3,4]、Wang 和 Rao^[5] 分别讨论了基于核实数据下线性 EV 模型, 部分线性 EV 模型以及半参数条件密度函数中的参数估计问题。Stute et al^[6] 利用核实数据讨论了非线性 EV 模型的经验似然推断问题。

本文通过利用核实数据, 构造了未知参数的两种经验对数似然比统计量, 即估计的经验对数似然比统计量和调整的经验对数似然比统计量。证明了所构造的经验似然比统计量渐近于 χ^2 分布, 所得结果可以用来构造未知参数的置信域。

收稿日期: 2008-03-10. 作者简介: 刘强(1976年9月生), 男, 博士, 副教授. 研究方向: 非参数统计与金融数据分析.

*基金项目: 北京市属高等学校人才强教计划资助项目; 首都经济贸易大学科研重点立项资助项目; 国家自然科学基金(10871013); 国家教育部博士点专项基金(20070005003).

2 方法与主要结果

假定 $E(e|\tilde{X}, X) = 0$ 。由于 $g(\cdot)$ 未知, 要估计 β_0 , 首先需要给出 $g(\cdot)$ 的估计。假定主要数据 $\{(Y_i, \tilde{X}_i)_{i=1}^N\}$ 是 N 个 i.i.d. 样本, 且与 i.i.d. 的核数据 $\{(Y_j, \tilde{X}_j, X_j)_{j=N+1}^{N+n}\}$ 相互独立。为了模型 (1) 的可识别性, 可以采用 Xue 和 Zhu^[1] 中“去一分量”的方法。记 $\beta \triangleq (\beta_1, \dots, \beta_{p-1})^T$, 不妨假定 β_0 的第 p 个分量 β_p 为正, 否则令 $\beta_p = -(1 - \|\beta\|^2)^{1/2}$ 。于是

$$\beta_0 = (\beta^T, (1 - \|\beta\|^2)^{1/2})^T.$$

因而对 β_0 的推断可以转化为对 $p-1$ 维未知参数 β 的推断。

首先, 利用局部线性光滑的方法给出未知函数 $g(\cdot)$ 和 $g'(\cdot)$ 的估计。对于任何 β , $g(\cdot)$ 和 $g'(\cdot)$ 的估计通过最小化下式给出。

$$\sum_{j=N+1}^{N+n} (Y_j - a - b(\beta_0^T X_j - t))^2 K_1\left(\frac{\beta_0^T X_j - t}{h_{1n}}\right), \quad (2)$$

其中 $K_1(\cdot)$ 为核函数, h_{1n} 为窗宽。记 \hat{a} 和 \hat{b} 为加权最小二乘问题 (2) 的解。未知函数 $g(\cdot)$ 和 $g'(\cdot)$ 的估计定义为 $\hat{g}(t; \beta) = \hat{a}$ 和 $\hat{g}'(t; \beta) = \hat{b}$, 经过简单计算有

$$\hat{g}(t; \beta) = \sum_{i=N+1}^{N+n} W_{ni}(t; \beta) Y_i, \quad \hat{g}'(t; \beta) = \sum_{i=N+1}^{N+n} \tilde{W}_{ni}(t; \beta) Y_i,$$

这里

$$W_{ni}(t; \beta) = U_{ni}(t; \beta, h_{1n}) / \sum_{j=N+1}^{N+n} U_{nj}(t; \beta, h_{1n}),$$

$$\tilde{W}_{ni}(t; \beta) = \tilde{U}_{ni}(t; \beta, h_{1n}) / \sum_{j=N+1}^{N+n} U_{nj}(t; \beta, h_{1n}),$$

$$U_{ni}(t; \beta, h_{1n}) = h_{1n}^{-1} K_1\left(\frac{\beta_0^T X_i - t}{h_{1n}}\right) (\beta_0^T X_i - t) [S_{n,2}(t; \beta, h_{1n}) - (\beta_0^T X_i - t) S_{n,1}(t; \beta, h_{1n})],$$

$$\tilde{U}_{ni}(t; \beta, h_{1n}) = h_{1n}^{-1} K_1\left(\frac{\beta_0^T X_i - t}{h_{1n}}\right) (\beta_0^T X_i - t) [(\beta_0^T X_i - t) S_{n,0}(t; \beta, h_{1n}) - S_{n,1}(t; \beta, h_{1n})],$$

$$S_{n,l}(t; \beta, h_{1n}) = n^{-1} h_{1n}^{-1} \sum_{i=N+1}^{N+n} (\beta_0^T X_i - t)^l K_1\left(\frac{\beta_0^T X_i - t}{h_{1n}}\right), \quad l = 0, 1, 2.$$

为了使用主要数据中的替代变量 \tilde{X} , 使得 Y 和 \tilde{X} 相关, 我们将模型 (1) 等价地转换为

$$Y = m(\tilde{X}, \beta) + \eta, \quad (3)$$

其中

$$m(\tilde{x}, \beta) = E[g(\beta_0^T X) | \tilde{X} = \tilde{x}], \quad \eta = e + g(\beta_0^T X) - E[g(\beta_0^T X) | \tilde{X}].$$

记 $J_\beta = \partial \beta_0 / \partial \beta$, 从而

$$m^{(1)}(\tilde{x}, \beta) \triangleq \frac{\partial}{\partial \beta} m(\tilde{x}, \beta) = E[J_\beta X g'(\beta_0^T X) | \tilde{X} = \tilde{x}].$$

构造辅助随机向量

$$Z_i(\beta) = (Y_i - m(\tilde{X}_i, \beta))m^{(1)}(\tilde{X}_i, \beta),$$

注意到 $EZ_i(\beta) = 0$ 当且仅当 β 为真参数, 从而检验 β 是否为真实参数的问题就转化为检验是否有等式 $EZ_i(\beta) = 0$ 对任意的 $i = 1, 2, \dots, N$ 成立。因此, 利用信息 $EZ_i(\beta) = 0$, 我们提出对数经验似然比函数

$$l(\beta) \triangleq l_n(\beta) = -2 \max \left\{ \sum_{i=1}^N \log(Np_i) \mid p_i \geq 0, \sum_{i=1}^N p_i = 1, \sum_{i=1}^N p_i Z_i(\beta) = 0 \right\}. \quad (4)$$

可以证明, 当 β 为真参数时, $l(\beta)$ 的分布渐近自由度为 $p-1$ 的标准卡方分布 χ_{p-1}^2 。然而 $l(\beta)$ 不能直接用于 β 的推断问题, 原因是 $l(\beta)$ 中含有未知函数 $m(\tilde{X}, \beta)$, $m^{(1)}(\tilde{X}, \beta)$ 。为此我们需要给出两个未知函数的估计。一个自然的方法是利用高维核估计给出上述函数的估计, 然而在非参数回归估计中, 这一推断需要较大的核实数据才能达到合适的精度。受王启华和 Härdle^[7] 的启发, 因而我们提出了如下降维方法。假设

$$m(\tilde{X}, \beta) = \mu(\alpha^T \tilde{X}, \beta),$$

其中 $\mu(\cdot)$ 为未知函数, α 为未知参数。这一假设的合理性可以参见王启华和 Härdle^[7]。从而 $m(\tilde{x}, \beta)$ 的估计可以定义为

$$\frac{\sum_{j=N+1}^{N+n} K_2\left(\frac{\alpha^T(\tilde{X}_j - \tilde{x})}{h_{2n}}\right) g(\beta_0^T X_j)}{\sum_{j=N+1}^{N+n} K_2\left(\frac{\alpha^T(\tilde{X}_j - \tilde{x})}{h_{2n}}\right)},$$

这里 $K_2(\cdot)$ 为核函数, h_{2n} 为收敛于 0 的窗宽。然而上述估计中的 α 和 $g(\cdot)$ 未知, 不能直接用来估计函数 $m(\tilde{x}, \beta)$ 。为解决这个问题, 一个自然的办法就是利用 α 和 $g(\cdot)$ 的相合估计来替代。

注意到 $E[g(\beta_0^T X) | \tilde{X}] = E[Y | \tilde{X}]$, 由 Zhu 和 Fang^[8], 可以利用切片逆回归核方法给出 α 的估计。定义

$$\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)^T, \quad R(Y) = (R_1(Y), \dots, R_p(Y))^T = (E[\tilde{X}_1 | Y], \dots, E[\tilde{X}_p | Y])^T,$$

为符号简洁, 我们也用 $\tilde{X}_i, i = N+1, \dots, N+n$ 表示 \tilde{X} 的 n 个独立复制, \tilde{X}_{ij} 为 \tilde{X}_i 的第 j 个分量, $j = 1, \dots, p$ 。记

$$\tilde{J}_{jn}(y) = \frac{1}{nh_{2n}} \sum_{i=N+1}^{N+n} \tilde{X}_{ij} K_3\left(\frac{y - Y_i}{h_{3n}}\right), \quad j = 1, \dots, p,$$

$$\hat{f}_{n,Y}(y) = \frac{1}{nh_{3n}} \sum_{i=N+1}^{N+n} K_3\left(\frac{y - Y_i}{h_{3n}}\right),$$

这里 $K_3(\cdot)$ 为核函数, h_{3n} 为窗宽。设 $\{\nu_n\}$ 为某个正的常数列, 记

$$\hat{f}_{bn,Y}(y) = \max\{\hat{f}_{n,Y}(y), \nu_n\}, \quad \hat{J}_{jn}(y) = \tilde{J}_{jn}(y) / \hat{f}_{bn,Y}(y),$$

$$\hat{\Lambda}_n = \frac{1}{n} \sum_{j=N+1}^{N+n} \hat{J}_{jn}(Y_j) \hat{J}_{jn}(Y_j)^T - \left(\frac{1}{n} \sum_{j=N+1}^{N+n} \hat{J}_{jn}(Y_j) \right) \left(\frac{1}{n} \sum_{j=N+1}^{N+n} \hat{J}_{jn}(Y_j) \right)^T.$$

设 α_n 为对应于 $\hat{\Lambda}_n$ 的最大特征值的特征根, 从而 α_n 可以作为 α 的相合估计. 记

$$\hat{R}_n(\tilde{x}, \beta) = \frac{1}{nh_{2n}} \sum_{j=N+1}^{N+n} \hat{g}(\beta_0^T X_j) K_2\left(\frac{\alpha_n^T (\tilde{X}_j - \tilde{x})}{h_{2n}}\right),$$

$$\hat{f}_n(\tilde{x}) = \frac{1}{nh_{2n}} \sum_{j=N+1}^{N+n} K_2\left(\frac{\alpha_n^T (\tilde{X}_j - \tilde{x})}{h_{2n}}\right).$$

记 $\hat{f}_{bn}(\tilde{x}) = \max\{\hat{f}_n(\tilde{x}), b_n\}$, 其中 $\{b_n\}$ 为某个正的常数列. 从而 $m(\tilde{x}, \beta)$ 的截尾估计为

$$\hat{m}(\tilde{x}, \beta) = \frac{\hat{R}_n(\tilde{x}, \beta)}{\hat{f}_{bn}(\tilde{x})}. \quad (5)$$

将 (5) 式中含有的 $\hat{g}(\beta_0^T X_j)$ 换为 $J_\beta X_j \hat{g}'(\beta_0^T X_j)$ 可以得到 $m^{(1)}(\tilde{x}, \beta)$ 的估计 $\hat{m}^{(1)}(\tilde{x}, \beta)$.

若记

$$\hat{Z}_i(\beta) = (Y_i - \hat{m}(\tilde{X}_i, \beta)) \hat{m}^{(1)}(\tilde{X}_i, \beta),$$

这里 $\hat{m}^{(1)}(\tilde{X}_i, \beta)$ 为 $\hat{m}(\tilde{X}_i, \beta)$ 对 β 的偏导数, 则定义估计的经验对数似然比为

$$\hat{l}(\beta) \triangleq \hat{l}_n(\beta) = -2 \max \left\{ \sum_{i=1}^N \log(Np_i) \mid p_i \geq 0, \sum_{i=1}^N p_i = 1, \sum_{i=1}^N p_i \hat{Z}_i(\beta) = 0 \right\}. \quad (6)$$

利用 Lagrange 乘数法, p_i 的最优值为 $p_i = N^{-1}(1 + \lambda^T \hat{Z}_i(\beta))^{-1}$, 其中 λ 为下述方程的解.

$$N^{-1} \sum_{i=1}^N \frac{\hat{Z}_i(\beta)}{1 + \lambda^T \hat{Z}_i(\beta)} = 0. \quad (7)$$

从而

$$\hat{l}(\beta) = 2 \sum_{i=1}^N \log(1 + \lambda^T \hat{Z}_i(\beta)). \quad (8)$$

设

$$V_0(\beta) = E[(Y - m(\tilde{X}, \beta))^2 m^{(1)}(\tilde{X}, \beta) m^{(1)T}(\tilde{X}, \beta)],$$

$$V(\beta) = V_0(\beta) + \gamma E[(m(\tilde{X}, \beta) - g(\beta_0^T X))^2 m^{(1)}(\tilde{X}, \beta) m^{(1)T}(\tilde{X}, \beta)],$$

其中 $\gamma = \lim \frac{N}{n}$. 为了记号简洁, 我们用 $m^{(1)T}(\tilde{X}, \beta)$ 表示 $m^{(1)}(\tilde{X}, \beta)$ 的转置. 与标准的经验似然对数函数不同, 由于 $\hat{Z}_i(\beta)$ ($i = 1, 2, \dots, n$) 不是 i.i.d., 因而 $\hat{l}(\beta)$ 的渐近分布不是标准 χ_{p-1}^2 分布.

定理 2.1 若满足第三节中的条件 C, 当 β 为真参数时, 则有

$$\hat{l}(\beta) \xrightarrow{L} \omega_1 \chi_{1,1}^2 + \omega_2 \chi_{1,2}^2 + \dots + \omega_{p-1} \chi_{1,p-1}^2,$$

其中 “ \xrightarrow{L} ” 表示依分布收敛, $\omega_1, \omega_2, \dots, \omega_{p-1}$ 为 $V_0^{-1}(\beta)V(\beta)$ 的特征值, $\chi_{1,1}^2, \chi_{1,2}^2, \dots, \chi_{1,p-1}^2$ 为相互独立的自由度为 1 的标准 χ^2 随机变量.

为了应用定理 2.1 构造 β 的置信域, 需要给出未知权重 ω_i 的相合估计。记

$$\begin{aligned}\hat{V}_0(\beta) &= N^{-1} \sum_{i=1}^N (Y_i - \hat{m}(\tilde{X}_i, \beta))^2 \hat{m}^{(1)}(\tilde{X}_i, \beta) \hat{m}^{(1)T}(\tilde{X}_i, \beta), \\ \hat{V}(\beta) &= \hat{V}_0(\beta) + \frac{N}{n^2} \sum_{j=N+1}^{N+n} (\hat{m}(\tilde{X}_j, \beta) - \hat{g}(\beta_0^T X_j))^2 \hat{m}^{(1)}(\tilde{X}_j, \beta) \hat{m}^{(1)T}(\tilde{X}_j, \beta).\end{aligned}$$

由此可以推出 $\hat{V}_0^{-1}(\hat{\beta})V(\hat{\beta})$ 的特征值 $\hat{\omega}_i$ 可以作为 ω_i 的相合估计, 其中 $\hat{\beta}$ 为极大化 $\hat{l}(\beta)$ 而得到的极大经验似然估计。用 $H(\cdot)$ 表示给定 $\{(Y_i, \tilde{X}_i)_{i=1}^N\}$ 和 $\{(Y_j, \tilde{X}_j, X_j)_{j=N+1}^{N+n}\}$ 下 $\omega_1 \chi_{1,1}^2 + \omega_2 \chi_{1,2}^2 + \cdots + \omega_{p-1} \chi_{1,p-1}^2$ 的条件分布, \hat{c}_α 为 $H(\cdot)$ 的 $1 - \alpha$ 分位点, 则 β 的渐近覆盖概率为 $1 - \alpha$ 的置信域为

$$\hat{I}_\alpha(\beta) = \{\beta : \hat{l}(\beta) \leq \hat{c}_\alpha\}. \quad (9)$$

事实上, 可以使用 Monte Carlo 模拟方法从 χ_1^2 分布中分别产生独立样本 $\chi_{1,1}^2, \chi_{1,2}^2, \cdots, \chi_{1,p-1}^2$ 而得到条件分布 $H(\cdot)$ 。

利用定理 2.1 构造 β 的置信域需要估计权重 ω_i , 这样会降低置信域的精度。按照 Rao 和 Scott^[9] 的结果, 可以推出

$$\tilde{\rho}(\beta) \sum_{i=1}^{p-1} \omega_i \chi_{1,i}^2$$

的分布渐近于一个自由度为 $p-1$ 的 χ^2 分布, 这里 $\tilde{\rho}(\beta) = (p-1)/\text{tr}\{\hat{V}_0^{-1}(\beta)\hat{V}(\beta)\}$ 。然而这种逼近仍然依赖于 ω_i , 下面给出一个调整的经验对数似然比。注意到

$$\tilde{\rho}(\beta) = \frac{\text{tr}\{\hat{V}^{-1}(\beta)\hat{V}(\beta)\}}{\text{tr}\{\hat{V}_0^{-1}(\beta)\hat{V}(\beta)\}},$$

检查 $\hat{l}(\beta)$ 的渐近表达式, 在 $\tilde{\rho}(\beta)$ 中用

$$\hat{B}(\beta) = \left\{ \sum_{i=1}^N \hat{Z}_i(\beta) \right\} \left\{ \sum_{i=1}^N \hat{Z}_i(\beta) \right\}^T,$$

替代 $\hat{V}(\beta)$, 从而得到一个新的调整因子 $\hat{\rho}(\beta)$ 。由此定义调整的对数经验似然比 (AEL)

$$\hat{l}_{ad}(\beta) = \hat{\rho}(\beta)\hat{l}(\beta). \quad (10)$$

定理 2.2 若满足第三节中的条件 C, 当 β 为真参数时, 则 $\hat{l}_{ad}(\beta) \xrightarrow{L} \chi_{p-1}^2$ 。

根据定理 2.2, 可以构造 β 的 $1 - \alpha$ 置信域

$$\hat{I}_{ad,\alpha}(\beta) = \{\beta : \hat{l}_{ad}(\beta) \leq \hat{c}_\alpha\}. \quad (11)$$

3 定理的证明

设 c 表示不依赖于 n 和 N 的正的常数, 在不同的地方可以表示不同的值。设 \mathcal{M}^2 是 \mathbf{R}^2 上或其子域上具有连续二阶偏导数组成的类。对任何的 p 维向量 a , 设 a_s 表示 a 的第 s 个分量, $s = 1, 2, \cdots, p$ 。

为了得到 $\hat{l}(\beta)$ 的渐近分布, 需要如下条件。

条件C 下列条件统称为条件C。

C1 $\beta_0^T X$ 的密度函数 $f(t)$ 在 \mathcal{T} 上满足 1 阶 Lipchitz 条件, 这里 $\mathcal{T} = \{t = \beta_0^T x : x \in A\}$, 其中 A 为 X 的有界支撑集合;

$$0 < \inf_{t \in \mathcal{T}} f(t) \leq \sup_{t \in \mathcal{T}} f(t) < \infty.$$

C2 (I) $g(t) \in \mathcal{M}^2$.

$$(II) \sup_{\tilde{x}} E[g^2(\beta_0 X) | \tilde{X} = \tilde{x}] < \infty, \sup_{\tilde{x}} E[g'(\beta_0 X)^2 | \tilde{X} = \tilde{x}] < \infty.$$

C3 (I) 核 $K_1(u)$ 是一个有界且对称的概率密度函数, 且满足

$$\int u^2 K_1(u) du \neq 0, \quad \int u^8 K_1(u) du < \infty.$$

(II) $K_2(u)$ 为可导的具有有界支撑的 2 阶核函数。

(III) $K_3(\cdot)$ 为对称的具有有界支撑 $[-1, 1]$ 的核函数, 且满足

$$\int u K_1(u) du = 0, \quad \int u^i K_1(u) du = 0, \quad i = 1, 2, 3.$$

C4 (I) $m(\tilde{x}, \beta) \in \mathcal{M}^2$, $m_s^{(1)}(\tilde{x}, \beta) \in \mathcal{M}^2$, $s = 1, 2, \dots, p-1$.

$$(II) \sup_{\tilde{x}} E[m^2(\tilde{x}, \beta)] < \infty, \sup_{\tilde{x}} E[m_s^{(1)}(\tilde{x}, \beta)^2] < \infty.$$

C5 \tilde{X} 的密度函数 $f_{\tilde{X}}(\tilde{x}) \in \mathcal{M}^2$, 且

$$\limsup_{n \rightarrow \infty} nP\{f_{\tilde{X}}(\tilde{x}) < b_n\} < \infty.$$

C6 (I) $\sup_{\tilde{x}} E(e^2 | \tilde{X} = \tilde{x}) < \infty$;

$$(II) E(\tilde{X}_s^4) < \infty; \sup_{\tilde{x}} E(X_s^4 | \tilde{X} = \tilde{x}) < \infty; \sup_{\tilde{x}} E(Y^2 | \tilde{X} = \tilde{x}) < \infty.$$

C7 (I) $nh_{1n}^3 h_{2n} b_n^2 \rightarrow \infty$, $h_{1n}^2 b_n^{-2} h_{2n}^{-1} \rightarrow 0$, $nh_{2n}^2 b_n^4 \rightarrow \infty$, $nh_{2n}^4 b_n^{-2} \rightarrow 0$.

(II) 对满足 $1/8 + c_2/4 < c_1 < 1/4 - c_2$ 的正数 c_1, c_2 和 $n \rightarrow \infty$, 有 $h_{3n} \sim n^{-c_1}$, $\nu_n \sim n^{-c_2}$, 这里 \sim 表示等价的两个量。

C8 记 $h_i(y) = R_i(y)f_Y(y)$, $i = 1, \dots, p$. $h_i(y)$, $f_Y(y)$ 3 阶可微且满足: 存在原点的某个领域 U_1 和 $c > 0$, 对于任意的 $u \in U_1$ 有

$$|f_Y^{(3)}(y+u) - f_Y^{(3)}(y)| \leq c|u|; \quad |h_i^{(3)}(y+u) - h_i^{(3)}(y)| \leq c|u|.$$

对任意的 $1 \leq i, j \leq p$ 和 $u \in U_1$ 有

$$|R_i(y+u)R_j(y+u) - R_i(y)R_j(y)| \leq c|u|.$$

C9 对任意的 $1 \leq i, j \leq p$ 有, $\sqrt{n}ER_i(Y)R_j(Y)I(f_Y(Y) < \nu_n) = o(1)$, 其中 $I(\cdot)$ 为示性函数。

C10 $V_0(\beta)$ 为一正定矩阵。

C11 $\gamma = \lim \frac{N}{n}$. 其中 $\gamma > 0$ 为常数。

注: 条件 **C1-C3** 是为了获取 $\hat{g} - g = O_p(h_{1n}^2 + (nh_{1n})^{-\frac{1}{2}})$, 条件 **C1** 主要用于 $\hat{g}(\cdot)$ 和 $\hat{g}'(\cdot)$ 的分母以高的概率偏离 0, 条件 **C2** 的使用原因是在估计 $g(\cdot)$ 时使用了 2 阶的核。条件 **C3** 保证 $\hat{g}(\cdot)$ 和 $\hat{g}'(\cdot)$ 有快的收敛速度, 具体见文献 [1,10]; 条件 **C4-C5**、**C7(I)** 是处理核实数据问题的基本条件, 条件 **C5** 主要是对 \tilde{X} 的密度函数 $f_{\tilde{X}}(\tilde{x})$ 的尾巴进行控制, 否则会给理论证明带来一定的困难, 见文献 [3,4] 和文献 [6]; 条件 **C6(II)**、**C7(II)** 和 **C8-C9** 是为了获取 $\hat{\alpha}_n - \alpha = O_p(n^{-\frac{1}{2}})$, 条件 **C8** 主要用于给出 Y 的密度函数的光滑度, 条件 **C7(II)** 和 **C8-C9** 已被王启华和 Härdle^[7] 与 Zhu 和 Fang^[8] 使用。其余条件是标准的。

我们把定理证明的主要步骤分成下面几个引理。

引理 1^[10] 若条件 **C1-C3** 成立, 如果 $h_{1n} = cn^{-\alpha}$, $0 < \alpha < 1/2$, $c > 0$, 则对于任何的正整数 $r \geq 2$, $N+1 \leq i \leq N+n$, 有

$$E \left[\left| \sum_{j=N+1}^{N+n} W_{nj}(\beta_0^T X_i; \beta) g(\beta_0^T X_j) - g(\beta_0^T X_i) \right|^r \right] = O(h_{1n}^{2r}),$$

$$E \left[\left| \sum_{j=N+1}^{N+n} \tilde{W}_{nj}(\beta_0^T X_i; \beta) g(\beta_0^T X_j) - g'(\beta_0^T X_i) \right|^r \right] = O(h_{1n}^r).$$

引理 2^[10] 在引理 1 的条件下, 对于任何的正整数 $r \geq 2$, $N+1 \leq i \leq N+n$, 有

$$E \left[|W_{ni}(\beta_0^T X_i; \beta)|^r \right] = O((nh_{1n})^{-r}),$$

$$E \left[\sum_{j=N+1, j \neq i}^{N+n} |W_{nj}(\beta_0^T X_i; \beta)|^r \right] = O((nh_{1n})^{1-r}),$$

$$E \left[|\tilde{W}_{ni}(\beta_0^T X_i; \beta)|^r \right] = O((nh_{1n})^{-r}) + O((n^3 h_{1n}^5)^{-r/2}),$$

$$E \left[\sum_{j=N+1, j \neq i}^{N+n} |\tilde{W}_{nj}(\beta_0^T X_i; \beta)|^r \right] = O(n^{1-r} h_{1n}^{1-2r}).$$

引理 3 在引理 1 的条件下, 对任意的 $N+1 \leq i \leq N+n$, 有

$$E \left[\hat{g}(\beta_0^T X_i; \beta) - g(\beta_0^T X_i) \right]^2 = O(h_{1n}^4) + O((nh_{1n})^{-1}), \quad (12)$$

$$E \left[\hat{g}'(\beta_0^T X_i; \beta) - g'(\beta_0^T X_i) \right]^2 = O(h_{1n}^2) + O((nh_{1n}^3)^{-1}). \quad (13)$$

证明 根据引理1和引理2, 有

$$\begin{aligned} & E [\hat{g}(\beta_0^T X_i; \beta) - g(\beta_0^T X_i)]^2 \\ &= E \left[\sum_{j=N+1}^{N+n} W_{nj}(\beta_0^T X_i; \beta) g(\beta_0^T X_j) - g(\beta_0^T X_i) \right]^2 + E \left[\sum_{j=N+1}^{N+n} W_{nj}(\beta_0^T X_i; \beta) e_j \right]^2 \\ &\leq ch_{1n}^4 + c \sum_{j=N+1}^{N+n} E [W_{nj}^2(\beta_0^T X_i; \beta)] \leq ch_{1n}^4 + c(nh_{1n})^{-1}. \end{aligned}$$

(12)式得证, 类似可以证明(13)式。

引理4 若条件C成立, 对任意的 $1 \leq i \leq N$, 有

$$\begin{aligned} \text{(i)} \quad & E [(\hat{m}(\tilde{X}_i, \beta) - m(\tilde{X}_i, \beta))^2 | \tilde{X}_i] \\ &\leq c(nh_{2n}b_n^2)^{-1} + ch_{2n}^4b_n^{-2} + cI[f_{\tilde{X}}(\tilde{X}_i) < 2b_n] + c(h_{1n}^4b_n^{-2}h_{2n}^{-1}) + c(nh_{1n}b_n^2h_{2n})^{-1}, \\ \text{(ii)} \quad & E [(\hat{m}_s^{(1)}(\tilde{X}_i, \beta) - m_s^{(1)}(\tilde{X}_i, \beta))^2 | \tilde{X}_i] \\ &\leq c(nh_{2n}b_n^2)^{-1} + ch_{2n}^4b_n^{-2} + cI[f_{\tilde{X}}(\tilde{X}_i) < 2b_n] + c(h_{1n}^2b_n^{-2}h_{2n}^{-1}) + c(nh_{1n}^3b_n^2h_{2n})^{-1}. \end{aligned}$$

证明 (i)与(ii)证明类似, 仅给出(i)的证明。设

$$f_{bn}(\cdot) = \max(f_{\tilde{X}}(\cdot), b_n), \quad m_{bn}(\cdot) = \frac{m(\cdot)f_{\tilde{X}}(\cdot)}{f_{bn}(\cdot)},$$

$\tilde{m}(\tilde{x}, \beta)$ 是用 α 替代 $\hat{m}(\tilde{x}, \beta)$ 的 α_n 得到的, 则有

$$\begin{aligned} & E [(m_{bn}(\tilde{X}_i, \beta) - m(\tilde{X}_i, \beta))^2 | \tilde{X}_i] \\ &= E \left[m^2(\tilde{X}_i, \beta) \left(1 - \frac{f_{\tilde{X}}(\tilde{X}_i)}{f_{bn}(\tilde{X}_i)} \right)^2 | \tilde{X}_i \right] \leq cI\{f_{\tilde{X}}(\tilde{X}_i) < b_n\}. \end{aligned}$$

下面仅需证明

$$\begin{aligned} & E [(\hat{m}(\tilde{X}_i, \beta) - m_{bn}(\tilde{X}_i, \beta))^2 | \tilde{X}_i] \\ &\leq c(nh_{2n}b_n^2)^{-1} + ch_{2n}^4b_n^{-2} + cI[f_{\tilde{X}}(\tilde{X}_i) < 2b_n] + c(h_{1n}^4b_n^{-2}h_{2n}^{-1}) + c(nh_{1n}b_n^2h_{2n})^{-1}, \quad (14) \end{aligned}$$

对任意的 $i = 1, 2, \dots, N$ 成立即可。记

$$\begin{aligned} J_{1n}(\tilde{x}) &= \frac{1}{nh_{2n}} \sum_{j=N+1}^{N+n} [g(\beta_0^T X_j) - m(\tilde{X}_j, \beta)] K_2 \left(\frac{\alpha^T (\tilde{X}_j - \tilde{x})}{h_{2n}} \right), \\ J_{2n}(\tilde{x}) &= \frac{1}{nh_{2n}} \sum_{j=N+1}^{N+n} [m(\tilde{X}_j, \beta) - m(\tilde{x}, \beta)] K_2 \left(\frac{\alpha^T (\tilde{X}_j - \tilde{x})}{h_{2n}} \right), \\ J_{3n}(\tilde{x}) &= \frac{1}{nh_{2n}} \sum_{j=N+1}^{N+n} [\hat{g}(\beta_0^T X_j) - g(\beta_0^T X_j)] K_2 \left(\frac{\alpha^T (\tilde{X}_j - \tilde{x})}{h_{2n}} \right), \\ J_{4n}(\tilde{x}) &= [\hat{f}_n(\tilde{x})f_{bn}(\tilde{x}) - \hat{f}_{bn}(\tilde{x})f_{\tilde{X}}(\tilde{x})]m(\tilde{x}, \beta) / [\hat{f}_{bn}(\tilde{x})f_{bn}(\tilde{x})]. \end{aligned}$$

记 $T_n(\tilde{x}) = \hat{m}(\tilde{x}, \beta) - m_{bn}(\tilde{x}, \beta)$, 则有

$$T_n(\tilde{x}) = [J_{1n}(\tilde{x}) + J_{2n}(\tilde{x}) + J_{3n}(\tilde{x})]/\hat{f}_{bn}(\tilde{x}) + J_{4n}(\tilde{x}) + [\hat{m}(\tilde{x}, \beta) - \tilde{m}(\tilde{x}, \beta)].$$

从而对于任意的 $i = 1, 2, \dots, N$, 有

$$\begin{aligned} E[J_{1n}^2(\tilde{X}_i) | \tilde{X}_i] &\leq 4b_n^{-2} E[J_{1n}^2(\tilde{X}_i) | \tilde{X}_i] + 4b_n^{-2} E[J_{2n}^2(\tilde{X}_i) | \tilde{X}_i] \\ &\quad + 4b_n^{-2} E[J_{3n}^2(\tilde{X}_i) | \tilde{X}_i] + 4E[J_{4n}^2(\tilde{X}_i) | \tilde{X}_i] \\ &\quad + E[(\hat{m}(\tilde{X}_i, \beta) - \tilde{m}(\tilde{X}_i, \beta))^2 | \tilde{X}_i]. \end{aligned} \quad (15)$$

经过简单计算, 对于任意的 $i = 1, 2, \dots, N$, 有

$$E[J_{1n}^2(\tilde{X}_i) | \tilde{X}_i] \leq c(nh_{2n}^p)^{-1}, \quad E[J_{2n}^2(\tilde{X}_i) | \tilde{X}_i] \leq c(nh_{2n}^p)^{-1} + ch_{2n}^{2k}. \quad (16)$$

$$E[J_{3n}^2(\tilde{X}_i) | \tilde{X}_i] \leq ch_{1n}^4 h_{2n}^{-p} b_n^{-2} + c(nh_{1n} h_{2n}^p b_n^2)^{-1}. \quad (17)$$

$$E[J_{4n}^2(\tilde{X}_i) | \tilde{X}_i] \leq cI\{f_{\tilde{X}}(\tilde{X}_i) < 2b_n\} + c(nh_{2n}^p b_n^2)^{-1} + ch_{2n}^{2k} b_n^{-2}. \quad (18)$$

由 Zhu 和 Fang^[8] 可知, $\alpha_n - \alpha = O_p(n^{-1/2})$, 经过一系列复杂运算, 可以证明

$$E[(\hat{m}(\tilde{X}_i, \beta) - \tilde{m}(\tilde{X}_i, \beta))^2 | \tilde{X}_i] \leq c(nh_{2n} b_n^2)^{-1} + ch_{2n}^4 b_n^{-2}.$$

结合 (15)-(18) 即可证明 (14) 式。

引理 5 在定理 2.1 的条件下, 若 β 为真实参数, 则有

$$N^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \xrightarrow{L} N(0, V(\beta)).$$

证明 经过一系列复杂运算, 可以证明

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \hat{Z}_i(\beta) &= n^{-\frac{1}{2}} \sum_{j=N+1}^{N+n} [m(\tilde{X}_j, \beta) - g(X_j, \beta)] m^{(1)}(\tilde{X}_j, \beta) \\ &\quad + N^{-\frac{1}{2}} \sum_{i=1}^N m^{(1)}(\tilde{X}_i, \beta) [Y_i - m(\tilde{X}_i, \beta)] + o_p(1). \end{aligned}$$

从而引理 5 得证。

引理 6 定理 2.1 的条件下, 若 β 为真实参数, 则有

$$\frac{1}{N} \sum_{i=1}^N \hat{Z}_i(\beta) \hat{Z}_i^T(\beta) \xrightarrow{P} V_0(\beta),$$

其中 $V_0(\beta)$ 的定义见定理 2.1。

证明 记

$$R_n(Y_i, \tilde{X}_i; \beta) = [m(\tilde{X}_i, \beta) - \hat{m}(\tilde{X}_i, \beta)] m^{(1)}(\tilde{X}_i, \beta) + Y_i [\hat{m}^{(1)}(\tilde{X}_i, \beta) - m^{(1)}(\tilde{X}_i, \beta)],$$

则有

$$\begin{aligned}
 & \frac{1}{N} \sum_{i=1}^N \hat{Z}_i(\beta) \hat{Z}_i^T(\beta) \\
 &= \frac{1}{N} \sum_{i=1}^N \eta_i^2 m^{(1)}(\tilde{X}_i, \beta) m^{(1)T}(\tilde{X}_i, \beta) + \frac{1}{N} \sum_{i=1}^N R_n(Y_i, \tilde{X}_i; \beta) R_n^T(Y_i, \tilde{X}_i; \beta) \\
 & \quad + \frac{1}{N} \sum_{i=1}^N R_n(Y_i, \tilde{X}_i; \beta) \eta_i m^{(1)T}(\tilde{X}_i, \beta) + \frac{1}{N} \sum_{i=1}^N \eta_i m^{(1)}(\tilde{X}_i, \beta) R_n^T(Y_i, \tilde{X}_i; \beta) \\
 & \triangleq \sum_{i=1}^4 \Lambda_i.
 \end{aligned}$$

利用大数定律得 $\Lambda_1 \xrightarrow{P} V_0(\beta)$ 。由引理4及标准的讨论可以证明 $\Lambda_i \xrightarrow{P} 0$, $i = 2, 3, 4$ 。从而引理6得证。

引理7 在定理2.1的条件下, 若 β 为真实参数, 则有

$$\max_{1 \leq i \leq N} \|\hat{Z}_i(\beta)\| = o_p(N^{\frac{1}{2}}).$$

证明 记 $\hat{Z}_{is}(\beta)$ 为 $\hat{Z}_i(\beta)$ 的第 s ($s = 1, 2, \dots, p-1$) 个分量,

$$\begin{aligned}
 \max_{1 \leq i \leq N} |\hat{Z}_{is}(\beta)| &\leq \max_{1 \leq i \leq N} |\eta_i m_s^{(1)}(\tilde{X}_i, \beta)| + \max_{1 \leq i \leq N} |[m(\tilde{X}_i, \beta) - \hat{m}(\tilde{X}_i, \beta)] m_s^{(1)}(\tilde{X}_i, \beta)| \\
 &\quad + \max_{1 \leq i \leq N} |[m_s^{(1)}(\tilde{X}_i, \beta) - \hat{m}_s^{(1)}(\tilde{X}_i, \beta)] \eta_i| \\
 &\quad + \max_{1 \leq i \leq N} |[m(\tilde{X}_i, \beta) - \hat{m}(\tilde{X}_i, \beta)] [m_s^{(1)}(\tilde{X}_i, \beta) - \hat{m}_s^{(1)}(\tilde{X}_i, \beta)]| \\
 &\triangleq \sum_{i=1}^4 B_i.
 \end{aligned}$$

由文献[11]中引理3可知, 若 $EY^2 < \infty$, 则

$$\max_{1 \leq i \leq n} |Y_i|/\sqrt{n} \xrightarrow{a.s.} 0.$$

由题设可知 $E[\eta_i m_s^{(1)}(\tilde{X}_i, \beta)]^2 < \infty$, 从而有 $B_1 = o_p(N^{\frac{1}{2}})$ 。利用Markov不等式, 结合引理4, 对于任意的 $\delta > 0$, 可以证明

$$\begin{aligned}
 P(N^{-1/2} B_2 > \delta) &\leq \sum_{i=1}^N P(|[m(\tilde{X}_i, \beta) - \hat{m}(\tilde{X}_i, \beta)] m_s^{(1)}(\tilde{X}_i, \beta)| > \delta \sqrt{N}) \\
 &\leq \frac{1}{N \delta^2} \sum_{i=1}^N E|[m(\tilde{X}_i, \beta) - \hat{m}(\tilde{X}_i, \beta)] m_s^{(1)}(\tilde{X}_i, \beta)|^2 \rightarrow 0.
 \end{aligned}$$

从而有 $B_2 = o_p(N^{\frac{1}{2}})$ 。类似地可以证明

$$B_3 = o_p(N^{\frac{1}{2}}), \quad B_4 = o_p(N^{\frac{1}{2}}).$$

引理 7 得证。

引理 8 在定理 2.1 的条件下, 若 β 为真参数, 有 $\lambda = O_p(N^{-\frac{1}{2}})$ 。

证明 由引理 5, 可以推出

$$N^{-1} \sum_{i=1}^N \hat{Z}_i(\beta) = O_p(n^{-\frac{1}{2}}).$$

由上式和引理 6, 利用 Owen^[11] 中 (1.14) 式相同的论证方法可以证明引理 8。

定理 2.1 的证明 在 (8) 式中利用 Taylor 展开, 并结合引理 6-8, 有

$$\hat{l}(\beta) = 2 \sum_{i=1}^N \left[\lambda^T \hat{Z}_i(\beta) - \frac{1}{2} \left(\lambda^T \hat{Z}_i(\beta) \right)^2 \right] + o_p(1). \quad (19)$$

又因为

$$0 = \sum_{i=1}^N \frac{\hat{Z}_i(\beta)}{1 + \lambda^T \hat{Z}_i(\beta)} = \sum_{i=1}^N \hat{Z}_i(\beta) - \sum_{i=1}^N \hat{Z}_i(\beta) \hat{Z}_i^T(\beta) + \sum_{i=1}^N \frac{\hat{Z}_i(\beta) (\lambda^T \hat{Z}_i(\beta))^2}{1 + \lambda^T \hat{Z}_i(\beta)}.$$

由引理 6-8 可知

$$\sum_{i=1}^N \frac{(\lambda^T \hat{Z}_i(\beta))^3}{1 + \lambda^T \hat{Z}_i(\beta)} = o_p(1).$$

从而

$$\begin{aligned} \sum_{i=1}^N \lambda^T \hat{Z}_i(\beta) &= \sum_{i=1}^N (\lambda^T \hat{Z}_i(\beta))^2 + o_p(1), \\ \lambda &= \left(\sum_{i=1}^N \hat{Z}_i(\beta) \hat{Z}_i^T(\beta) \right)^{-1} \sum_{i=1}^N \hat{Z}_i(\beta) + o_p(N^{-\frac{1}{2}}). \end{aligned}$$

结合 (19) 式, 可以证明

$$\hat{l}(\beta) = \left(N^{-\frac{1}{2}} V^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \right)^T V^{\frac{1}{2}} V_0^{-1} V^{\frac{1}{2}} \left(N^{-\frac{1}{2}} V^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \right) + o_p(1).$$

根据引理 6 可知

$$N^{-\frac{1}{2}} V^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \xrightarrow{L} N(0, I_p). \quad (20)$$

而 $V^{\frac{1}{2}} V_0^{-1} V^{\frac{1}{2}}$ 和 $V_0^{-1} V$ 具有相同的特征值, 从而定理 1 得证。

定理 2.2 的证明 回顾 $\hat{l}_{ad}(\beta)$ 的定义, 结合 (19) 式可推得

$$\hat{l}_{ad}(\beta) = \left(N^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \right)^T \hat{V}(\beta) \left(N^{-\frac{1}{2}} \sum_{i=1}^N \hat{Z}_i(\beta) \right) + o_p(1).$$

类似于引理 6 的证明, 可以证明 $\hat{V}(\beta) \xrightarrow{P} V(\beta)$, 结合 (20) 式可以证明 $\hat{l}_{ad}(\beta)$ 依分布收敛于自由度为 $p-1$ 的 χ^2 分布。

参考文献:

- [1] Xue L G, Zhu L X. Empirical likelihood for single-index models[J]. *Journal of Multivariate Analysis*, 2006, 97: 1295-1312
- [2] Sepanski J H, Lee L F. Semiparametric estimation of nonlinear error-in-variables models with validation study[J]. *Journal of Nonparametric Statistics*, 1995, 4: 365-394
- [3] Wang Q H. Estimation of partial linear error-in-variables models with validation data[J]. *Journal of Multivariate Analysis*, 1999, 69: 30-64
- [4] Wang Q H. Likelihood-based kernel estimation in semiparametric errors-in-variables models with validation data[J]. *Journal of Multivariate Analysis*, 2007, 98: 455-480
- [5] Wang Q H, Rao J N K. Empirical likelihood-based inference in linear errors-in-covariables models with validation data[J]. *Biometrika*, 2002, 89(2): 345-358
- [6] Stute W, et al. Empirical likelihood inference in nonlinear errors-in-covariables models with validation data[J]. *Journal of the American Statistical Association*, 2007, 102(477): 332-346
- [7] 王启华, Härdle W. 核实数据帮助下误差在反映线性模型经验似然降维推断[J]. *中国科学, A 辑*, 2004, 34(5): 549-566
Wang Q H, Härdle W. Empirical likelihood-based dimension reduction inference for linear error-in-responses models with validation study[J]. *Science in China Series A*, 2004, 47(6): 921-939
- [8] Zhu L X, Fang K T. Asymptotics for kernel estimator of sliced inverse regression[J]. *The Annals of Statistics*, 1996, 24: 1053-1068
- [9] Rao J N K, Scott A J. The analysis of categorical data from complex sample surveys: chi-squared tests for goodness of fit and independence in two-way tables[J]. *Journal of the American Statistical Association*, 1981, 76: 221-230
- [10] Zhu L X, Xue L G. Empirical likelihood confidence regions in a partially linear single-index model[J]. *Journal of the American Statistical Association*, 2006, 68: 549-570
- [11] Owen A B. Empirical likelihood ratio confidence regions[J]. *The Annals of Statistics*, 1990, 18(1): 90-120

Empirical Likelihood-based Inference for Single-index EV Models with Validation Data

LIU Qiang¹, XUE Liu-gen²

(1- School of Statistics, Capital University of Economics and Business, Beijing 100070;

2- College of Applied Sciences, Beijing University of Technology, Beijing 100124)

Abstract: The single-index EV model is considered in this paper. When the explanatory variable is erroneously measured, especially the relationship between the surrogate variable and the explanatory one is unknown, the curse of dimension often appears in the estimation in the single-index model. In this paper, we propose a dimension reduction method, by which the above problem can be settled effectively. The empirical likelihood approach is used to construct the confidence regions of the parameter. Because a semiparametric dimension reduction technique is employed in the estimation procedure, the results obtained in the paper can be used to deal with the dimension curse problem in the data processing in the fields such as biological engineering and network engineering, and has higher value of application.

Keywords: single-index EV models; empirical likelihood; validation data; dimension reduction

Received: 10 Mar 2008. **Accepted:** 12 Nov 2008.

Foundation item: The Funding Project for Academic Human Resources Development in Institutions of Higher Education Under the Jurisdiction of Beijing Municipality; the Science Research Foundation of Capital University of Economics and Business; the National Natural Science Foundation of China (10871013); the P.H. D. Program Foundation of Ministry of Education of China (20070005003).